

# Large Margin Multiclass Gaussian Classification with Differential Privacy

Manas A. Pathak and Bhiksha Raj

Carnegie Mellon University  
`{manasp, bhiksha}@cs.cmu.edu`

**Abstract.** As increasing amounts of sensitive personal information is aggregated into data repositories, it has become important to develop mechanisms for processing the data without revealing information about individual data instances. The differential privacy model provides a framework for the development and theoretical analysis of such mechanisms. In this paper, we propose an algorithm for learning a discriminatively trained multi-class Gaussian classifier that satisfies differential privacy using a large margin loss function with a perturbed regularization term. We present a theoretical upper bound on the excess risk of the classifier introduced by the perturbation.

## 1 Introduction

In recent years, vast amounts of personal data is being aggregated in the form of medical, financial records, social networks, and government census data. As these often contain sensitive information, a database curator interested in releasing a function such as a statistic evaluated over the data is faced with the prospect that it may lead to a breach of privacy of the individuals who contributed to the database. It is therefore important to develop techniques for retrieving desired information from a dataset without revealing any information about individual data instances. *Differential privacy* [1] is a theoretical model proposed to address this issue. A query mechanism evaluated over a dataset is said to satisfy differential privacy if it is likely to produce the same output on a dataset differing by at most one element. This implies that an adversary having complete knowledge of all data instances but one along with *a priori* information about the remaining instance, is not likely to be able to infer any more information about the remaining instance by observing the output of the mechanism.

One of the most common applications for such large data sets such as the ones mentioned above is for training classifiers that can be used to categorize new data. If the training data contains private data instances, an adversary should not be able to learn anything about the individual training dataset instances by analyzing the output of the classifier. Recently, mechanisms for learning differentially private classifiers have been proposed for logistic regression [2]. In this method, the objective function which is minimized by the classification algorithm is modified by adding a linear perturbation term. Compared to the original classifier, there is an additional error introduced by the perturbation term in the

differentially private classifier. It is important to have an upper bound on this error as a cost of preserving privacy.

The work mentioned above is largely restricted to binary classification, while multi-class classifiers are more useful in many practical situations. In this paper, we propose an algorithm for learning multi-class Gaussian classifiers which satisfies differential privacy. Gaussian classifiers that model the distributions of individual classes as being generated from Gaussian distribution or a mixture of Gaussian distributions [3] are commonly used as multi-class classifiers. We use a large margin discriminative algorithm for training the classifier introduced by Sha and Saul [4]. To ensure that the learned multi-class classifier preserves differential privacy, we modify the objective function by introducing a perturbed regularization term.

## 2 Differential Privacy

In recent years, the differential privacy model proposed by Dwork, *et al.* [1] has emerged as a robust standard for data privacy. It originated from the statistical database model, where the dataset  $D$  is a collection of elements and a randomized *query mechanism*  $M$  produces a response when performed on a given dataset. Two datasets  $D$  and  $D'$  differing by at most one element are said to be *adjacent*. There are two proposed definitions for adjacent datasets one based on symmetric difference –  $D'$  containing of one entry less than  $D$ , and one based on substitution – one entry of  $D'$  differs in value from  $D$ . We use the substitution definition of adjacency previously used by [5,2], where the one entry of the dataset  $D = \{x_1, \dots, x_{n-1}, x_n\}$  is modified to result in an adjacent dataset  $D' = \{x_1, \dots, x_{n-1}, x'_n\}$ . The query mechanism  $M$  is said to satisfy differential privacy if the probability of  $M$  resulting in a solution  $S$  when performed on a dataset  $D$  is very close to the probability of  $M$  resulting in the same solution  $S$  when executed on an adjacent dataset  $D'$ . Assuming the query mechanism to be a function  $M : D \mapsto \text{range}(M)$  with a probability function  $P$  defined over the space of  $M$ , differential privacy is formally defined as follows.

**Definition 1.** A randomized function  $M$  satisfies  $\epsilon$ -differential privacy if for all adjacent datasets  $D$  and  $D'$  and for any  $S \in \text{range}(M)$ ,

$$\left| \log \frac{P(M(D) = S)}{P(M(D') = S)} \right| \leq \epsilon.$$

The value of the  $\epsilon$  parameter, which is referred to as *leakage*, determines the degree of privacy. As there is always a trade-off between privacy and utility, the choice of  $\epsilon$  is motivated by the requirements of the application.

In a machine learning setting, the query mechanism can be thought of as an algorithm learning the classification, regression or density estimation rule which is evaluated over the training dataset. The output of an algorithm satisfying differential privacy is likely to be same when the value of any single dataset instance is modified, and therefore, no additional information can be obtained about any

individual training data instances with certainty by observing the output of the learning algorithm, beyond what is already known to an adversary. Differential privacy is a strong definition of privacy – it provides *ad omnia* guarantee as opposed to most other models that provide *ad hoc* guarantees against specific set of attacks and adversarial behaviors.

## 2.1 Related Work

The earlier work on differential privacy was related to functional approximations for simple data mining tasks and data release mechanisms [6,7,8,9]. Although many of these works have connection to machine learning problems, more recently the design and analysis of machine learning algorithms satisfying differential privacy has been actively studied. Kasiviswanathan, *et al.* [5] present a framework for converting a general agnostic PAC learning algorithm to an algorithm that satisfies privacy constraints. Chaudhuri and Monteleoni [2] use the exponential mechanism [10] to create a differentially private logistic regression classifier by adding Laplace noise to the estimated parameters. They propose another differentially private formulation which involves modifying the objective function of the logistic regression classifier by adding a linear term scaled by Laplace noise. The second formulation is advantageous because it is independent of the classifier sensitivity which difficult to compute in general and it can be shown that using a perturbed objective function introduces a lower error as compared to the exponential mechanism.

However, the above mentioned differentially private classification algorithms only address the problem of binary classification. Although it is possible to extend binary classification algorithms to multi-class using techniques like one-vs-all, it is much more expensive to do so as compared to a naturally multi-class classification algorithm. Jagannathan, *et al.* [11] present a differentially private random decision tree learning algorithm which can be applied to multi-class classification. Their approach involves perturbing leaf nodes using the sensitivity method, and they do not provide theoretical analysis of excess risk of the perturbed classifier. In this paper, we propose a modification to the naturally multi-class large margin Gaussian classification algorithm [4,12].

## 3 Large Margin Gaussian Classifiers

We investigate the large margin multi-class classification algorithm introduced by Sha and Saul [4]. The training dataset  $(\mathbf{x}, \mathbf{y})^1$  contains  $n$  iid  $d$ -dimensional training data instances  $\mathbf{x}_i \in \mathbb{R}^d$  each with labels  $y_i \in \{1, \dots, C\}$ . We consider the setting where each class is modeled as a single Gaussian ellipsoid. Each class ellipsoid is parametrized by the centroid  $\boldsymbol{\mu}_c \in \mathbb{R}^d$ , the inverse covariance matrix  $\boldsymbol{\Psi}_c \in \mathbb{R}^{d \times d}$ , and a scalar offset  $\theta_c \geq 0$ . The decision rule is to assign an instance

---

<sup>1</sup> Notation: vectors and matrices are denoted by **boldface**.

$\mathbf{x}_i$  to the class having smallest Mahalanobis distance [13] with the scalar offset from  $\mathbf{x}_i$  to the centroid of that class.

$$y_i = \operatorname{argmin}_c (\mathbf{x}_i - \boldsymbol{\mu}_c)^T \boldsymbol{\Psi}_c (\mathbf{x}_i - \boldsymbol{\mu}_c) + \theta_c. \quad (1)$$

To simplify the notation, we expand  $(\mathbf{x}_i - \boldsymbol{\mu}_c)^T \boldsymbol{\Psi}_c (\mathbf{x}_i - \boldsymbol{\mu}_c)$  and collect the parameters for each class as the following  $(d+1) \times (d+1)$  positive semidefinite matrix

$$\boldsymbol{\Phi}_c = \begin{bmatrix} \boldsymbol{\Psi}_c & -\boldsymbol{\Psi}_c \boldsymbol{\mu}_c \\ -\boldsymbol{\mu}_c^T \boldsymbol{\Psi}_c & \boldsymbol{\mu}_c^T \boldsymbol{\Psi}_c \boldsymbol{\mu}_c + \theta_c \end{bmatrix} \quad (2)$$

and also append a unit element to each  $d$ -dimensional vector  $\mathbf{x}_i$ . The decision rule for a data instance  $\mathbf{x}_i$  simplifies to

$$y_i = \operatorname{argmin}_c \mathbf{x}_i^T \boldsymbol{\Phi}_c \mathbf{x}_i. \quad (3)$$

The discriminative training procedure involves estimating a set of positive semidefinite matrices  $\{\boldsymbol{\Phi}_1, \dots, \boldsymbol{\Phi}_C\}$  from the training data  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  which optimize the performance on the decision rule mentioned above. We apply the large margin intuition that the optimal classifier must maximize the distance of training data instances from the decision boundaries. This leads to the classification algorithm being robust to outliers with provably strong generalization guarantees. Formally, we require that for each training data instance  $\mathbf{x}_i$  with label  $y_i$ , the distance from  $\mathbf{x}_i$  to the centroid of class  $y_i$  is at least less than its distance from centroids of all other classes by one.

$$\forall c \neq y_i : \mathbf{x}_i^T \boldsymbol{\Phi}_c \mathbf{x}_i \geq 1 + \mathbf{x}_i^T \boldsymbol{\Phi}_{y_i} \mathbf{x}_i.$$

Analogous to support vector machines, the training algorithm is an optimization problem minimizing the *hinge loss* denoted by  $[f]_+ = \max(0, f)$ , with a linear penalty for incorrect classification. We use the sum of traces of inverse covariance matrices for each classes as a *regularization* term. The regularization requires that if we can learn a classifier which labels every training data instance correctly, we choose the one with the lowest inverse covariance or highest covariance for each class ellipsoid as this prevents the classifier from over-fitting. The parameter  $\lambda$  controls the trade off between the loss function and the regularization.

$$J(\boldsymbol{\Phi}, \mathbf{x}, \mathbf{y}) = \sum_i \sum_{c \neq y_i} [1 + \mathbf{x}_i^T (\boldsymbol{\Phi}_{y_i} - \boldsymbol{\Phi}_c) \mathbf{x}_i]_+ + \lambda \sum_c \operatorname{trace}(\boldsymbol{\Psi}_c). \quad (4)$$

The inverse covariance matrix  $\boldsymbol{\Psi}_c$  is contained in the upper left size  $d \times d$  block of the matrix  $\boldsymbol{\Phi}_c$ . We replace it with  $\mathbf{I}_{\Phi} \boldsymbol{\Phi}_c \mathbf{I}_{\Phi}$ , where  $\mathbf{I}_{\Phi}$  is the truncated size  $(d+1) \times (d+1)$  identity matrix with the last diagonal element  $I_{\Phi_{d+1,d+1}}$  set to zero. The optimization problem becomes

$$\begin{aligned} J(\boldsymbol{\Phi}, \mathbf{x}, \mathbf{y}) &= \sum_i \sum_{c \neq y_i} [1 + \mathbf{x}_i^T (\boldsymbol{\Phi}_{y_i} - \boldsymbol{\Phi}_c) \mathbf{x}_i]_+ + \lambda \sum_c \operatorname{trace}(\mathbf{I}_{\Phi} \boldsymbol{\Phi}_c \mathbf{I}_{\Phi}) \\ &= L(\boldsymbol{\Phi}, \mathbf{x}, \mathbf{y}) + N(\boldsymbol{\Phi}). \end{aligned} \quad (5)$$

The hinge loss being non-differentiable is not very convenient for our analysis; we replace it with a surrogate loss function called Huber loss  $l_h$  [14] which has similar characteristics to the hinge loss for small values of  $h$ .

$$\ell_h(\Phi_c, \mathbf{x}_i, y_i) = \begin{cases} 0 & \text{if } \mathbf{x}_i^T (\Phi_c - \Phi_{y_i}) \mathbf{x}_i > h, \\ \frac{1}{4h} [h - \mathbf{x}_i^T (\Phi_{y_i} - \Phi_c) \mathbf{x}_i]^2 & \text{if } |\mathbf{x}_i^T (\Phi_c - \Phi_{y_i}) \mathbf{x}_i| \leq h \\ -\mathbf{x}_i^T (\Phi_{y_i} - \Phi_c) \mathbf{x}_i & \text{if } \mathbf{x}_i^T (\Phi_c - \Phi_{y_i}) \mathbf{x}_i < -h. \end{cases} \quad (6)$$

The objective function is convex function of positive semidefinite matrices  $\Phi_c$ . The optimization can be formulated as a semidefinite programming problem [15] and be solved efficiently using interior point methods.

The large margin classification framework can be easily extended to modeling each class with a mixture of Gaussians. Similar to support vector machines, when training with non-separable data, we can introduce slack parameters to permit margin violations. These extensions do not change the basic characteristics of the learning algorithm. The optimization problem remains to be a convex semidefinite program with piecewise linear terms and is equally tractable. For simplicity, we restrict our discussion to single Gaussians and hard margins in this paper. As we shall see, it is easy to extend our proposed modifications to these cases.

## 4 Differentially Private Large Margin Gaussian Classifiers

We modify the large margin Gaussian classification formulation to satisfy differential privacy by introducing a perturbation term in the objective function. As we will see in Section 5.2, this modification leads to a classifier that preserves differential privacy.

We generate the size  $(d+1) \times (d+1)$  perturbation matrix  $\mathbf{b}$  with density

$$P(\mathbf{b}) \propto \exp\left(-\frac{\epsilon}{2}\|\mathbf{b}\|\right), \quad (7)$$

where  $\|\cdot\|$  is the Frobenius norm (element-wise  $\ell_2$  norm) and  $\epsilon$  is the privacy parameter. One method of generating such a  $\mathbf{b}$  matrix is to sample the norm  $\|\mathbf{b}\|$  from  $\Gamma((d+1)^2, \frac{2}{\epsilon})$  and the direction of  $\mathbf{b}$  at random.

Our proposed learning algorithm minimizes the following objective function  $J_p(\Phi, \mathbf{x}, \mathbf{y})$ , where the subscript  $p$  denotes privacy.

$$\begin{aligned} J_p(\Phi, \mathbf{x}, \mathbf{y}) &= L(\Phi, \mathbf{x}, \mathbf{y}) + \lambda \sum_c \text{trace}(\mathbf{I}_\Phi \Phi_c \mathbf{I}_\Phi) + \sum_c \sum_{ij} b_{ij} \Phi_{cij} \\ &= J(\Phi, \mathbf{x}, \mathbf{y}) + \sum_c \sum_{ij} b_{ij} \Phi_{cij}. \end{aligned} \quad (8)$$

As the dimensionality of the perturbation matrix  $\mathbf{b}$  is same as that of the classifier parameters  $\Phi_c$ , the parameter space of  $\Phi$  does not change after perturbation. In other words, given two datasets  $(\mathbf{x}, \mathbf{y})$  and  $(\mathbf{x}', \mathbf{y}')$ , if  $\Phi^p$  minimizes

$J_p(\Phi, \mathbf{x}, \mathbf{y})$ , it is always possible to have  $\Phi^P$  minimize  $J_p(\Phi, \mathbf{x}', \mathbf{y}')$ . This is a necessary condition for the classifier  $\Phi^P$  satisfying differential privacy.

Furthermore, as the perturbation term is convex and positive semidefinite, the perturbed objective function  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  has the same properties as the unperturbed objective function  $J(\Phi, \mathbf{x}, \mathbf{y})$ . Also, the perturbation does not introduce any additional computational cost as compared to the original algorithm.

## 5 Theoretical Analysis

### 5.1 Proof of Differential Privacy

In the following theorem, we prove that the classifier minimizing the perturbed optimization function  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  satisfies  $\epsilon$ -differential privacy. Given the dataset  $(\mathbf{x}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{n-1}, y_{n-1}), (\mathbf{x}_n, y_n)\}$ , the probability of learning the classifier  $\Phi^P$  is close to the probability of learning the same classifier  $\Phi^P$  given its adjacent dataset  $(\mathbf{x}', \mathbf{y}') = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{n-1}, y_{n-1}), (\mathbf{x}'_n, y'_n)\}$  differing wlog on the  $n^{\text{th}}$  instance. As we mentioned in the previous section, it is always possible to find such a classifier  $\Phi^P$  minimizing both  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  and  $J_p(\Phi, \mathbf{x}', \mathbf{y}')$  due to the perturbation matrix being in the same space as the optimization parameters.

Our proof requires a strictly convex perturbed objective function resulting in a unique solution  $\Phi^P$  minimizing it. This in turn requires that the loss function  $L(\Phi, \mathbf{x}, y)$  is strictly convex and differentiable, and the regularization term  $N(\Phi)$  is convex. These seemingly strong constraints are satisfied by many commonly used classification algorithms such as logistic regression, support vector machines, and our general perturbation technique can be extended to those algorithms. In our proposed algorithm, the Huber loss is by definition a differentiable function and the trace regularization term is convex and differentiable. Additionally, we require that the difference in the gradients of  $L(\Phi, \mathbf{x}, y)$  calculated over for two adjacent training datasets is bounded. We prove this property in Lemma 1 given in the appendix.

**Theorem 1.** *For any two adjacent training datasets  $(\mathbf{x}, \mathbf{y})$  and  $(\mathbf{x}', \mathbf{y}')$ , the classifier  $\Phi^P$  minimizing the perturbed objective function  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  satisfies differential privacy.*

$$\left| \log \frac{P(\Phi^P | \mathbf{x}, \mathbf{y})}{P(\Phi^P | \mathbf{x}', \mathbf{y}')} \right| \leq \epsilon',$$

where  $\epsilon' = \epsilon + k$  for a constant factor  $k = \log \left( 1 + \frac{2\alpha}{n\lambda} + \frac{\alpha^2}{n^2\lambda^2} \right)$  with a constant value of  $\alpha$ .

*Proof.* As  $J(\Phi, \mathbf{x}, \mathbf{y})$  is convex and differentiable, there is a unique solution  $\Phi^*$  that minimizes it. As the perturbation term  $\sum_c \sum_{ij} b_{ij} \Phi_{cij}$  is also convex and differentiable, the perturbed objective function  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  also has a unique

solution  $\Phi^p$  that minimizes it. Differentiating  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  wrt  $\Phi_c$ , we have

$$\frac{\partial}{\partial \Phi_c} J_p(\Phi, \mathbf{x}, \mathbf{y}) = \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}, \mathbf{y}) + \lambda \mathbf{I}_\Phi + \mathbf{b}. \quad (9)$$

Substituting the optimal  $\Phi_c^p$  in the derivative gives us

$$\lambda \mathbf{I}_\Phi + \mathbf{b} = -\frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}, \mathbf{y}).$$

This relation shows that two different values of  $\mathbf{b}$  cannot result in the same optimal  $\Phi^p$ . As the perturbed objective function  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  is also convex and differentiable, there is a bijective map between the perturbation  $\mathbf{b}$  and the unique  $\Phi^p$  minimizing  $J_p(\Phi, \mathbf{x}, \mathbf{y})$ .

Let  $\mathbf{b}_1$  and  $\mathbf{b}_2$  be the two perturbations applied when training with the adjacent datasets  $(\mathbf{x}, \mathbf{y})$  and  $(\mathbf{x}', \mathbf{y}')$ , respectively. Assuming that we obtain the same optimal solution  $\Phi^p$  while minimizing both  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  with perturbation  $\mathbf{b}_1$  and  $J_p(\Phi, \mathbf{x}, \mathbf{y})$  with perturbation  $\mathbf{b}_2$ ,

$$\begin{aligned} \lambda \mathbf{I}_\Phi + \mathbf{b}_1 &= -\frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}, \mathbf{y}), \\ \lambda \mathbf{I}_\Phi + \mathbf{b}_2 &= -\frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}', \mathbf{y}'), \\ \mathbf{b}_1 - \mathbf{b}_2 &= \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}', \mathbf{y}') - \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}, \mathbf{y}). \end{aligned} \quad (10)$$

We apply Lemma 1 after taking Frobenius norm on both sides.

$$\begin{aligned} \|\mathbf{b}_1 - \mathbf{b}_2\| &= \left\| \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}', \mathbf{y}') - \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}, \mathbf{y}) \right\| \\ &= \left\| \sum_{i=1}^{n-1} \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}_i, y_i) + \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}'_n, y'_n) \right. \\ &\quad \left. - \sum_{i=1}^{n-1} \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}_i, y_i) - \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}_n, y_n) \right\| \\ &= \left\| \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}'_n, y'_n) - \frac{\partial}{\partial \Phi_c} L(\Phi^p, \mathbf{x}_n, y_n) \right\| \leq 2. \end{aligned}$$

Using this property, we can calculate the ratio of densities of drawing the perturbation matrices  $\mathbf{b}_1$  and  $\mathbf{b}_2$  as

$$\frac{P(\mathbf{b} = \mathbf{b}_1)}{P(\mathbf{b} = \mathbf{b}_2)} = \frac{\frac{1}{\text{surf}(\|\mathbf{b}_1\|)} \|\mathbf{b}_1\|^d \exp[-\frac{\epsilon}{2} \|\mathbf{b}_1\|]}{\frac{1}{\text{surf}(\|\mathbf{b}_2\|)} \|\mathbf{b}_2\|^d \exp[-\frac{\epsilon}{2} \|\mathbf{b}_2\|]},$$

where  $\text{surf}(\|\mathbf{b}\|)$  is the surface area of the  $(d+1)$ -dimensional hypersphere with radius  $\|\mathbf{b}\|$ . As  $\text{surf}(\|\mathbf{b}\|) = \text{surf}(1)\|\mathbf{b}\|^d$ , where  $\text{surf}(1)$  is the area of the unit

$(d + 1)$ -dimensional hypersphere, the ratio of the densities becomes

$$\frac{P(\mathbf{b} = \mathbf{b}_1)}{P(\mathbf{b} = \mathbf{b}_2)} = \exp\left[\frac{\epsilon}{2}(\|\mathbf{b}_2\| - \|\mathbf{b}_1\|)\right] \leq \exp\left[\frac{\epsilon}{2}\|\mathbf{b}_2 - \mathbf{b}_1\|\right] \leq \exp(\epsilon). \quad (11)$$

The ratio of the densities of learning  $\Phi^p$  using the adjacent datasets  $(\mathbf{x}, \mathbf{y})$  and  $(\mathbf{x}', \mathbf{y}')$  is given by

$$\frac{P(\Phi^p | \mathbf{x}, \mathbf{y})}{P(\Phi^p | \mathbf{x}', \mathbf{y}')} = \frac{P(\mathbf{b} = \mathbf{b}_1)}{P(\mathbf{b} = \mathbf{b}_2)} \frac{|\det(\mathbf{J}(\Phi^p \rightarrow \mathbf{b}_1 | \mathbf{x}, \mathbf{y}))|^{-1}}{|\det(\mathbf{J}(\Phi^p \rightarrow \mathbf{b}_2 | \mathbf{x}', \mathbf{y}'))|^{-1}}, \quad (12)$$

where  $\mathbf{J}(\Phi^p \rightarrow \mathbf{b}_1 | \mathbf{x}, \mathbf{y})$  and  $\mathbf{J}(\Phi^p \rightarrow \mathbf{b}_2 | \mathbf{x}', \mathbf{y}')$  are the Jacobian matrices of the bijective mappings from  $\Phi^p$  to  $\mathbf{b}_1$  and  $\mathbf{b}_2$ , respectively. Following a procedure identical to Theorem 2 of [16] (omitted due to lack of space), it can be shown that the ratio of Jacobian determinants is upper bounded by a constant factor  $\exp(k) = 1 + \frac{2\alpha}{n\lambda} + \frac{\alpha^2}{n^2\lambda^2}$  for a constant value of  $\alpha$ . Therefore, the ratio of the densities of learning  $\Phi^p$  using the adjacent datasets becomes

$$\frac{P(\Phi^p | \mathbf{x}, \mathbf{y})}{P(\Phi^p | \mathbf{x}', \mathbf{y}')} \leq \exp(\epsilon + k) = \exp(\epsilon'). \quad (13)$$

Similarly, we can show that the probability ratio is lower bounded by  $\exp(-\epsilon')$ , which together with Equation (13) satisfies the definition of differential privacy.

□

## 5.2 Analysis of Excess Error

In the remainder of this section, we denote the terms  $J(\Phi, \mathbf{x}, \mathbf{y})$  and  $L(\Phi, \mathbf{x}, \mathbf{y})$  by  $J(\Phi)$  and  $L(\Phi)$  respectively for conciseness. To establish a bound on excess risk of the classifier given by the proposed algorithm minimizing the perturbed objective function, in Lemma 2 we show that the objective function  $J(\Phi)$  satisfies strong convexity. The objective function  $J(\Phi)$  contains the loss function  $L(\Phi)$  computed over the training data  $(\mathbf{x}, \mathbf{y})$  and the regularization term  $N(\Phi)$  – this is known as the regularized *empirical risk* of the classifier  $\Phi$ . In the following theorem, we establish a bound on the regularized empirical excess risk of the differentially private classifier minimizing the perturbed objective function over the classifier minimizing the unperturbed objective function.

**Theorem 2.** *With probability at least  $1 - \delta$ , the regularized empirical excess risk of the classifier  $\Phi^p$  minimizing the perturbed objective function  $J_p(\Phi)$  over the classifier  $\Phi^*$  minimizing the unperturbed objective function  $J(\Phi)$  is bounded as*

$$J(\Phi^p) \leq J(\Phi^*) + \frac{8(d+1)^4 C}{\epsilon^2 \lambda} \log^2\left(\frac{d}{\delta}\right).$$

*Proof.* We use the definition of  $J_p(\Phi) = J(\Phi) + \sum_c \sum_{ij} b_{ij} \Phi_{cij}$  and the optimality of  $\Phi^p$ , i.e.,  $J_p(\Phi^p) \leq J_p(\Phi^*)$ .

$$\begin{aligned} J(\Phi^p) + \sum_c \sum_{ij} b_{ij} \Phi_{cij}^p &\leq J(\Phi^*) + \sum_c \sum_{ij} b_{ij} \Phi_{cij}^*, \\ J(\Phi^p) &\leq J(\Phi^*) + \sum_c \sum_{ij} b_{ij} (\Phi_{cij}^* - \Phi_{cij}^p). \end{aligned} \quad (14)$$

Using the strong convexity of  $J(\Phi)$  as given by Lemma 2 and the optimality of  $J(\Phi^*)$ , we have

$$\begin{aligned} J(\Phi^*) &\leq J\left(\frac{\Phi^p + \Phi^*}{2}\right) \leq \frac{J(\Phi^p) + J(\Phi^*)}{2} - \frac{\lambda}{8} \sum_c \|\Phi_c^* - \Phi_c^p\|^2, \\ J(\Phi^p) - J(\Phi^*) &\geq \frac{\lambda}{4} \sum_c \|\Phi_c^p - \Phi_c^*\|^2. \end{aligned} \quad (15)$$

Similarly, using the strong convexity of  $J_p(\Phi)$  and the optimality of  $J_p(\Phi^p)$ ,

$$\begin{aligned} J_p(\Phi^p) &\leq J_p\left(\frac{\Phi^p + \Phi^*}{2}\right) \leq \frac{J_p(\Phi^p) + J_p(\Phi^*)}{2} - \frac{\lambda}{8} \sum_c \|\Phi_c^p - \Phi_c^*\|^2, \\ J_p(\Phi^*) - J_p(\Phi^p) &\geq \frac{\lambda}{4} \sum_c \|\Phi_c^p - \Phi_c^*\|^2. \end{aligned}$$

Substituting the definition  $J_p(\Phi) = J(\Phi) + \sum_c \sum_{ij} b_{ij} \Phi_{cij}$ ,

$$\begin{aligned} J(\Phi^*) + \sum_c \sum_{ij} b_{ij} \Phi_{cij}^* - J(\Phi^p) - \sum_c \sum_{ij} b_{ij} \Phi_{cij}^p &\geq \frac{\lambda}{4} \sum_c \|\Phi_c^* - \Phi_c^p\|^2 \\ \sum_c \sum_{ij} b_{ij} (\Phi_{cij}^* - \Phi_{cij}^p) - (J(\Phi^p) - J(\Phi^*)) &\geq \frac{\lambda}{4} \sum_c \|\Phi_c^* - \Phi_c^p\|^2. \end{aligned}$$

Substituting the lower bound on  $J(\Phi^p) - J(\Phi^*)$  given by Equation (15),

$$\begin{aligned} \sum_c \sum_{ij} b_{ij} (\Phi_{cij}^* - \Phi_{cij}^p) &\geq \frac{\lambda}{2} \sum_c \|\Phi_c^* - \Phi_c^p\|^2, \\ \left[ \sum_c \sum_{ij} b_{ij} (\Phi_{cij}^* - \Phi_{cij}^p) \right]^2 &\geq \frac{\lambda^2}{4} \left[ \sum_c \|\Phi_c^* - \Phi_c^p\|^2 \right]^2. \end{aligned} \quad (16)$$

Using the Cauchy-Schwarz inequality, we have,

$$\left[ \sum_c \sum_{ij} b_{ij} (\Phi_{cij}^* - \Phi_{cij}^p) \right]^2 \leq C \|\mathbf{b}\|^2 \sum_c \|\Phi_c^* - \Phi_c^p\|^2 \quad (17)$$

Combining this with Equation (16) gives us

$$\begin{aligned} C\|\mathbf{b}\|^2 \sum_c \|\Phi_c^* - \Phi_c^p\|^2 &\geq \frac{\lambda^2}{4} \left[ \sum_c \|\Phi_c^* - \Phi_c^p\|^2 \right]^2, \\ \sum_c \|\Phi_c^* - \Phi_c^p\|^2 &\leq \frac{4C}{\lambda^2} \|\mathbf{b}\|^2. \end{aligned} \quad (18)$$

Combining this with Equation (17) gives us

$$\sum_c \sum_{ij} b_{ij} (\Phi_{cij}^* - \Phi_{cij}^p) \leq \frac{2C}{\lambda} \|\mathbf{b}\|^2.$$

We bound  $\|\mathbf{b}\|^2$  with probability at least  $1 - \delta$  as given by Lemma 4.

$$\sum_c \sum_{ij} b_{ij} (\Phi_{cij}^* - \Phi_{cij}^p) \leq \frac{8(d+1)^4 C}{\epsilon^2 \lambda} \log^2 \left( \frac{d}{\delta} \right). \quad (19)$$

Substituting this in Equation (14) proves the theorem.  $\square$

The upper bound on the regularized empirical risk is in  $O(\frac{C}{\epsilon^2})$ . The bound increases for smaller values of  $\epsilon$  which implies tighter privacy and therefore suggests a trade off between privacy and utility.

The regularized empirical risk of a classifier is calculated over a given training dataset. In practice, we are more interested in how the classifier will perform on new test data which is assumed to be generated from the same source as the training data. The expected value of the loss function computed over the data is called the *true risk*  $\tilde{L}(\Phi) = \mathbb{E}[L(\Phi)]$  of the classifier  $\Phi$ . In the following theorem, we establish a bound on the true excess risk of the differentially private classifier minimizing the perturbed objective function and the classifier minimizing the original objective function.

**Theorem 3.** *With probability at least  $1 - \delta$ , the true excess risk of the classifier  $\Phi^p$  minimizing the perturbed objective function  $J_p(\Phi)$  over the classifier  $\Phi^*$  minimizing the unperturbed objective function  $J(\Phi)$  is bounded as*

$$\begin{aligned} \tilde{L}(\Phi^p) &\leq \tilde{L}(\Phi^*) + \frac{4\sqrt{d}(d+1)^2 C}{\epsilon \lambda} \log \left( \frac{d}{\delta} \right) \\ &\quad + \frac{8(d+1)^4 C}{\epsilon^2 \lambda} \log^2 \left( \frac{d}{\delta} \right) + \frac{16}{\lambda n} \left[ 32 + \log \left( \frac{1}{\delta} \right) \right]. \end{aligned}$$

*Proof.* Let the expected value of the regularized empirical risk be

$$\tilde{J}(\Phi) = \tilde{L}(\Phi) + \lambda \sum_c \text{trace}(\mathbf{I}_\Phi \Phi_c \mathbf{I}_\Phi). \quad (20)$$

Let  $\Phi^r$  be the classifier minimizing  $\tilde{J}(\Phi)$ , i.e.,  $\tilde{J}(\Phi^r) \leq \tilde{J}(\Phi^*)$ .

Rearranging the terms, we have

$$\begin{aligned}\tilde{J}(\Phi^p) &= \tilde{J}(\Phi^*) + [\tilde{J}(\Phi^p) - \tilde{J}(\Phi^r)] + [\tilde{J}(\Phi^r) - \tilde{J}(\Phi^*)] \\ &\leq \tilde{J}(\Phi^*) + [\tilde{J}(\Phi^p) - \tilde{J}(\Phi^r)].\end{aligned}$$

Substituting the definition of  $\tilde{J}(\Phi)$ ,

$$\begin{aligned}\tilde{L}(\Phi^p) + \lambda \sum_c \text{trace}(\mathbf{I}_\Phi \Phi_c^p \mathbf{I}_\Phi) &\leq \tilde{L}(\Phi^*) + \lambda \sum_c \text{trace}(\mathbf{I}_\Phi \Phi_c^* \mathbf{I}_\Phi) + [\tilde{J}(\Phi^p) - \tilde{J}(\Phi^r)] \\ \tilde{L}(\Phi^p) &\leq \tilde{L}(\Phi^*) + \lambda \sum_c \text{trace}[\mathbf{I}_\Phi (\Phi_c^* - \Phi_c^p) \mathbf{I}_\Phi] + [\tilde{J}(\Phi^p) - \tilde{J}(\Phi^r)].\end{aligned}\quad (21)$$

From Lemma 3 and Equation (18), we have,

$$\begin{aligned}\left[ \sum_c \text{trace}[\mathbf{I}_\Phi (\Phi_c^* - \Phi_c^p) \mathbf{I}_\Phi] \right]^2 &\leq dC \sum_c \|\Phi_c - \Phi'_c\|^2 \\ &\leq \frac{4dC^2}{\lambda^2} \|\mathbf{b}\|^2 = \frac{16d(d+1)^4 C^2}{\epsilon^2 \lambda^2} \log^2 \left( \frac{d}{\delta} \right).\end{aligned}$$

Taking the square root,

$$\sum_c \text{trace}[\mathbf{I}_\Phi (\Phi_c^* - \Phi_c^p) \mathbf{I}_\Phi] \leq \frac{4\sqrt{d}(d+1)^2 C}{\epsilon \lambda} \log \left( \frac{d}{\delta} \right). \quad (22)$$

Sridharan, *et al.* [17] present a bound on the true excess risk of any classifier as an expression of the bound on the regularized empirical excess risk for that classifier. With probability at least  $1 - \delta$ ,

$$\tilde{J}(\Phi^p) - \tilde{J}(\Phi^r) \leq 2[J(\Phi^p) - J(\Phi^*)] + \frac{16}{\lambda n} \left[ 32 + \log \left( \frac{1}{\delta} \right) \right].$$

Substituting the bound from Theorem 2,

$$\tilde{J}(\Phi^p) - \tilde{J}(\Phi^r) \leq \frac{8(d+1)^4 C}{\epsilon^2 \lambda} \log^2 \left( \frac{d}{\delta} \right) + \frac{16}{\lambda n} \left[ 32 + \log \left( \frac{1}{\delta} \right) \right]. \quad (23)$$

Substituting the results from Equations (22) and (23) into Equation (21) proves the theorem.  $\square$

Similar to the bound on the regularized empirical excess risk, the bound on the true excess risk is also inversely proportional to  $\epsilon$  reflecting the privacy-utility trade-off. The bound is linear in the number of classes  $C$ , which is a consequence of the multi-class classification. The classifier learned using a higher value of the regularization parameter  $\lambda$  will have a higher covariance for each class ellipsoid. This would also make the classifier less sensitive to the perturbation. This intuition is confirmed by the fact that the true excess risk bound is inversely proportional to  $\lambda$ .

## 6 Conclusion

In this paper, we present a discriminatively trained Gaussian classification algorithm that satisfies differential privacy. Our proposed technique involves adding a perturbation term to the objective function. We prove that the proposed algorithm satisfies differential privacy and establish a bound on the excess risk of the classifier learned by the algorithm which is inversely proportional to the data dimensionality which is directly proportional to the number of classes and inversely proportional to the privacy parameter  $\epsilon$  reflecting a trade-off between privacy and utility.

In the future, we plan to extend this work along two main directions: extending our perturbation technique for a general class of learning algorithms and applying results from theory of large margin classifiers to arrive at tighter excess risk bounds for the differentially private large margin classifiers. Our intuition is that compared to other classification algorithms, a large margin classifier should be much more robust to perturbation. This would also give us insights into designing low error inducing mechanisms for differentially private classifiers.

## Acknowledgements

We would like to thank the anonymous reviewers for their insightful comments.

## References

1. Dwork, C.: Differential privacy. In: International Colloquium on Automata, Languages and Programming. (2006)
2. Chaudhuri, K., Monteleoni, C.: Privacy-preserving logistic regression. In: Neural Information Processing Systems. (2008) 289–296
3. McLachlan, G., Peel, D.: Finite Mixture Models. Wiley series in probability and statistics. Wiley-Interscience (2000)
4. Sha, F., Saul, L.K.: Large margin gaussian mixture modeling for phonetic classification and recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing. (2006) 265–268
5. Kasiviswanathan, S.P., Lee, H.K., Nissim, K., Raskhodnikova, S., Smith, A.: What can we learn privately? In: IEEE Symposium on Foundations of Computer Science. (2008) 531–540
6. Dinur, I., Nissim, K.: Revealing information while preserving privacy. In: Symposium on Principles of Database Systems. (2003)
7. Dwork, C., Nissim, K.: Privacy-preserving datamining on vertically partitioned databases. In: CRYPTO. (2004)
8. Blum, A., Dwork, C., McSherry, F., Nissim, K.: Practical privacy: The suLQ framework. In: Symposium on Principles of Database Systems. (2005)
9. Barak, B., Chaudhuri, K., Dwork, C., Kale, S., McSherry, F., Talwar, K.: Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In: Symposium on Principles of Database Systems. (2007) 273–282

10. Dwork, C., McSherry, F., Nissim, K., Smith, A.: Calibrating noise to sensitivity in private data analysis. In: Theory of Cryptography Conference. Volume 3876. (2006) 265–284
11. Jagannathan, G., Pillaipakkamnatt, K., Wright, R.N.: A practical differentially private random decision tree classifier. In: ICDM Workshop on Privacy Aspects of Data Mining. (2009) 114–121
12. Sha, F., Saul, L.K.: Large margin hidden markov models for automatic speech recognition. In: Neural Information Processing Systems. (2007) 1249–1256
13. Mahalanobis, P.C.: On the generalised distance in statistics. Proceedings of the National Institute of Sciences of India **2** (1936) 49–55
14. Chapelle, O.: Training a support vector machine in the primal. Neural Computation **19**(5) (2007) 1155–1178
15. Vandenberghe, L., Boyd, S.: Semidefinite programming. SIAM Review **38** (1996) 49–95
16. Chaudhuri, K., Monteleoni, C., Sarwate, A.D.: Differentially private empirical risk minimization. arXiv:0912.0071v4 [cs.LG] (2010)
17. Sridharan, K., Shalev-Shwartz, S., Srebro, N.: Fast rates for regularized objectives. In: Neural Information Processing Systems. (2008) 1545–1552

## Appendix

**Lemma 1.** *Assuming all the data instances to lie within a unit  $\ell_2$  ball, the difference in the derivative of Huber loss function  $L(\Phi, \mathbf{x}, y)$  calculated over two data instances  $(\mathbf{x}_i, y_i)$  and  $(\mathbf{x}'_i, y'_i)$  is bounded.*

$$\left\| \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}_i, y_i) - \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}'_i, y'_i) \right\| \leq 2.$$

*Proof.* The derivative of the Huber loss function for the data instance  $\mathbf{x}_i$  with label  $y_i$  is

$$\frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}_i, y_i) = \begin{cases} 0 & \text{if } \mathbf{x}_i^T (\Phi_c - \Phi_{y_i}) \mathbf{x}_i > h, \\ \frac{1}{2h} [h - \mathbf{x}_i^T (\Phi_{y_i} - \Phi_c) \mathbf{x}_i] \mathbf{x}_i \mathbf{x}_i^T & \text{if } |\mathbf{x}_i^T (\Phi_c - \Phi_{y_i}) \mathbf{x}_i| \leq h, \\ \mathbf{x}_i \mathbf{x}_i^T & \text{if } \mathbf{x}_i^T (\Phi_c - \Phi_{y_i}) \mathbf{x}_i < -h. \end{cases}$$

The data points lie in a  $\ell_2$  ball of radius 1,  $\forall i : \|\mathbf{x}_i\|_2 \leq 1$ . Using linear algebra, it is easy to show that the Frobenius norm of the matrix  $\mathbf{x}_i \mathbf{x}_i^T$  is same as the  $\ell_2$  norm of the vector  $\mathbf{x}_i$ ,  $\|\mathbf{x}_i \mathbf{x}_i^T\| = \|\mathbf{x}_i\|_2 \leq 1$ .

As the term  $\frac{1}{2h} [h - \mathbf{x}_i^T (\Phi_{y_i} - \Phi_c) \mathbf{x}_i]$  is at most one when  $|\mathbf{x}_i^T (\Phi_c - \Phi_{y_i}) \mathbf{x}_i| \leq h$ , the Frobenius norm of the derivative of the Huber loss function is at most one in all cases,  $\left\| \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}_i, y_i) \right\| \leq 1$ . Using a similar argument for data instance  $\mathbf{x}'_i$  with label  $y'_i$ , we have  $\left\| \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}'_i, y'_i) \right\| \leq 1$ .

Finally, using the triangle inequality  $\|\mathbf{a} - \mathbf{b}\| = \|\mathbf{a} + (-\mathbf{b})\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|$ ,

$$\begin{aligned} & \left\| \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}_i, y_i) - \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}'_i, y'_i) \right\| \\ & \leq \left\| \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}_i, y_i) \right\| + \left\| \frac{\partial}{\partial \Phi_c} L(\Phi, \mathbf{x}'_i, y'_i) \right\| \leq 2. \end{aligned}$$

□

**Lemma 2.** *The objective function  $J(\Phi)$  is  $\lambda$ -strongly convex. For  $0 \leq \alpha \leq 1$ ,*

$$J(\alpha\Phi + (1 - \alpha)\Phi') \leq \alpha J(\Phi) + (1 - \alpha)J(\Phi') - \frac{\lambda\alpha(1 - \alpha)}{2} \sum_c \|\Phi_c - \Phi'_c\|^2.$$

*Proof.* By definition, Huber loss is  $\lambda$ -strongly convex, i.e.

$$L(\alpha\Phi + (1 - \alpha)\Phi') \leq \alpha L(\Phi) + (1 - \alpha)L(\Phi') - \frac{\lambda\alpha(1 - \alpha)}{2} \|\Phi - \Phi'\|^2. \quad (24)$$

where the Frobenius norm of the matrix set  $\Phi - \Phi'$  is the sum of norms of the component matrices  $\Phi_c - \Phi'_c$ ,

$$\|\Phi - \Phi'\|^2 = \sum_c \|\Phi_c - \Phi'_c\|^2. \quad (25)$$

As the regularization term  $N(\Phi)$  is linear,

$$\begin{aligned} N(\alpha\Phi + (1 - \alpha)\Phi') &= \lambda \sum_c \text{trace}(\alpha \mathbf{I}_\Phi \Phi_c \mathbf{I}_\Phi + (1 - \alpha) \mathbf{I}_\Phi \Phi'_c \mathbf{I}_\Phi) \\ &= \alpha\lambda \sum_c \text{trace}(\mathbf{I}_\Phi \Phi_c \mathbf{I}_\Phi) + (1 - \alpha)\lambda \sum_c \text{trace}(\mathbf{I}_\Phi \Phi'_c \mathbf{I}_\Phi) \\ &= \alpha N(\Phi) + (1 - \alpha)N(\Phi'). \end{aligned} \quad (26)$$

The lemma follows directly from the definition  $J(\Phi) = L(\Phi) + N(\Phi)$ .

□

**Lemma 3.**

$$\frac{1}{dC} \left[ \sum_c \text{trace}[\mathbf{I}_\Phi (\Phi_c - \Phi'_c) \mathbf{I}_\Phi] \right]^2 \leq \sum_c \|\Phi_c - \Phi'_c\|^2$$

*Proof.* Let  $\Phi_{c,i,j}$  be the  $(i, j)^{\text{th}}$  element of the size  $(d+1) \times (d+1)$  matrix  $\Phi_c - \Phi'_c$ . By the definition of the Frobenius norm, and using the identity  $N \sum_{i=1}^N x_i^2 \geq (\sum_{i=1}^N x_i)^2$ ,

$$\begin{aligned} \sum_c \|\Phi_c - \Phi'_c\|^2 &= \sum_c \sum_{i=1}^{d+1} \sum_{j=1}^{d+1} \Phi_{c,i,j}^2 \geq \sum_c \sum_{i=1}^{d+1} \Phi_{c,i,i}^2 \geq \sum_c \sum_{i=1}^d \Phi_{c,i,i}^2 \\ &\geq \frac{1}{dC} \left( \sum_c \sum_{i=1}^d \Phi_{c,i,i} \right)^2 = \frac{1}{dC} \left[ \sum_c \text{trace}[\mathbf{I}_\Phi (\Phi_c - \Phi'_c) \mathbf{I}_\Phi] \right]^2. \end{aligned}$$

□

**Lemma 4.**

$$P \left[ \|\mathbf{b}\| \geq \frac{2(d+1)^2}{\epsilon} \log \left( \frac{d}{\delta} \right) \right] \leq \delta.$$

*Proof.* Similar to the union bound argument used in Lemma 5 in [2].